

# Terminology finding in the Sketch Engine: an Evaluation

**Adam Kilgarriff**

Lexical Computing Ltd.

## **ABSTRACT**

The [Sketch Engine](#) is a leading corpus query tool, in use for lexicography at OUP, CUP, Collins and Le Robert, and at national language institutes of eight countries, and for teaching and research in many universities. Its distinctive feature is the 'word sketch' a one page, automatic, corpus, derived summary of a word's grammatical and collocational behaviour.

Very large corpora and word sketches are available for sixty languages.

A number of tools and resources have recently been added with translators and terminologists in mind. The resources are parallel corpora: EUROPARL-7 and the various datasets available in the OPUS collection. The tools are bilingual word sketches and the term finder. These have been reported on in previous Asling/Aslib conferences.

One remaining task is to make the Sketch Engine functions conveniently accessible to translators and terminologists. We have recently done this via IntelliWebSearch, a tool which lets the user highlight text in the environment they are working in, which could be a CAT tool or Microsoft Word, and, with a key sequence, query their preferred database. So now the key sequence can take the translator or terminologist to a browser window showing the word sketch, or parallel concordance, or any of a number of other reports, for the expression they are working on.

## **1. The term-finding functionality**

The term-finder starts from a domain corpus, and a reference corpus. First it finds all the noun phrases, and their frequencies, in both corpora. It then takes the ratio, and the items with highest ratios will be terms, as in Figure 1 (data supplied by lead users, the World Intellectual Property Organisation).

Term	Frequency	Freq/mill	Score
station de base	<a href="#">28612</a>	3292.2	3293.2
station mobile	<a href="#">12514</a>	1439.9	1440.9
communication sans fil	<a href="#">8189</a>	942.3	943.3
liaison montante	<a href="#">6561</a>	754.9	737.5
terminal mobile	<a href="#">7406</a>	852.2	709.8
liaison descendante	<a href="#">5434</a>	625.3	626.3
stations de base	<a href="#">5010</a>	576.5	577.5
réseau de communication	<a href="#">4255</a>	489.6	490.6
communication mobile	<a href="#">4722</a>	543.3	462.5
point d' accès	<a href="#">3907</a>	449.6	450.6
modes de réalisation	<a href="#">3486</a>	401.1	402.1
réseau d' accès	<a href="#">3241</a>	372.9	373.9
réseau sans fil	<a href="#">2903</a>	334.0	335.0
accès radio	<a href="#">2412</a>	277.5	278.5
transfert intercellulaire	<a href="#">2408</a>	277.1	278.1

14

**Figure 1: French terms in the mobile communications domain**

In some cases, as with WIPO, the user will have domain corpora, but in others they will not. In that case they may use the BootCaT procedure (Baroni and Bernardini 2004). The user, typically a translator working in a domain where they are not an expert, inputs a few domain-specific 'seed words'; these are sent to a search engine, and the hits identified by the search engine are gathered, cleaned, de-duplicated and processed to give a domain-specific corpus. This functionality has been found to support translators well (Bernardini et al 2013). For some time, the Sketch Engine has incorporated a BootCaT tool, allowing users to create an instant corpus for a domain, which means they can then compare this corpus with a reference corpus to find the keywords of the domain. The functionality has recently been extended so the user can find the terms alongside key words. Thus, where the user has Bootcatted an English environment corpus, the Sketch Engine provides the "key words and terms" report shown in Figure 2.

The requirements for the term-finding functionality are:

- a processing chain, comprising tokeniser, lemmatiser and part-of-speech tagger, installed and ready to apply to the user's domain corpus
- a reference corpus processed with the processing chain
- a term grammar.

At time of writing, these are in place for Chinese, English, French, German, Japanese, Korean, Russian, Spanish and Portuguese.

Keywords		Terms
<input type="checkbox"/> dioxide (415.2, <a href="#">427</a> )	<input type="checkbox"/> mutualism (75.6, <a href="#">8</a> )	<input type="checkbox"/> carbon dioxide (567.1)
<input type="checkbox"/> trophic (264.9, <a href="#">33</a> )	<input type="checkbox"/> radiative (75.0, <a href="#">12</a> )	<input type="checkbox"/> greenhouse effect (515.0)
<input type="checkbox"/> greenhouse (238.4, <a href="#">282</a> )	<input type="checkbox"/> gasses (75.0, <a href="#">12</a> )	<input type="checkbox"/> water vapor (486.8)
<input type="checkbox"/> ecology (237.7, <a href="#">196</a> )	<input type="checkbox"/> lca (74.4, <a href="#">10</a> )	<input type="checkbox"/> global warming (298.8)
<input type="checkbox"/> methane (233.5, <a href="#">108</a> )	<input type="checkbox"/> biotic (74.2, <a href="#">10</a> )	<input type="checkbox"/> industrial ecology (261.6)
<input type="checkbox"/> arrhenius (232.2, <a href="#">25</a> )	<input type="checkbox"/> acidification (74.1, <a href="#">9</a> )	<input type="checkbox"/> infrared radiation (170.9)
<input type="checkbox"/> photosynthesis (230.6, <a href="#">46</a> )	<input type="checkbox"/> above-ground (73.6, <a href="#">9</a> )	<input type="checkbox"/> carbon cycle (169.0)
<input type="checkbox"/> callendar (215.4, <a href="#">22</a> )	<input type="checkbox"/> holism (73.5, <a href="#">9</a> )	<input type="checkbox"/> surface temperature (161.0)
<input type="checkbox"/> ecosystems (211.4, <a href="#">114</a> )	<input type="checkbox"/> felzer (73.5, <a href="#">7</a> )	<input type="checkbox"/> elevated carbon (156.4)
<input type="checkbox"/> warming (193.8, <a href="#">504</a> )	<input type="checkbox"/> carbonic (72.4, <a href="#">9</a> )	<input type="checkbox"/> elevated carbon dioxide (156.4)
<input type="checkbox"/> keeling (192.5, <a href="#">23</a> )	<input type="checkbox"/> loa (71.5, <a href="#">10</a> )	<input type="checkbox"/> greenhouse gas (135.8)
<input type="checkbox"/> carbon (186.8, <a href="#">558</a> )	<input type="checkbox"/> biogeography (71.2, <a href="#">9</a> )	<input type="checkbox"/> climate system (134.1)
<input type="checkbox"/> n't (177.1, <a href="#">17</a> )	<input type="checkbox"/> organisms (70.4, <a href="#">86</a> )	<input type="checkbox"/> food web (124.3)
<input type="checkbox"/> gases (173.9, <a href="#">159</a> )	<input type="checkbox"/> mauna (69.7, <a href="#">10</a> )	<input type="checkbox"/> amount of carbon dioxide (116.8)
<input type="checkbox"/> -oct- (169.3, <a href="#">28</a> )	<input type="checkbox"/> flowering (68.4, <a href="#">23</a> )	<input type="checkbox"/> other greenhouse (114.2)
<input type="checkbox"/> vapor (151.3, <a href="#">72</a> )	<input type="checkbox"/> emitted (68.2, <a href="#">27</a> )	<input type="checkbox"/> global temperature (109.1)
<input type="checkbox"/> deforestation (144.7, <a href="#">38</a> )	<input type="checkbox"/> suess (67.4, <a href="#">7</a> )	<input type="checkbox"/> atmospheric carbon (107.1)
<input type="checkbox"/> ecosystem (138.6, <a href="#">88</a> )	<input type="checkbox"/> infrared (65.1, <a href="#">44</a> )	<input type="checkbox"/> human activity (106.7)

**Figure 2: English keywords and terms in the environment domain. The tickboxes are so the user can iterate the procedure to extend and refine the corpus**

## 2. Evaluation

To evaluate a term-finder for a language and a domain, a list of all the 'true terms' is required. Then we can compute precision and recall.

One problem: how to define the domain? The straightforward answer: provide a corpus of it. Then we have the more constrained task of assessing recall and precision, from a given corpus, when the terms in that corpus are known.

Another problem: won't two different terminologists inevitably deliver two different lists?

We approached the task by hunting for research datasets comprising domain corpora and term-lists derived, by human experts, from them. In most cases, this had been done as part of a term-finding task, so there were also published papers, with term-finding results, over these datasets, so we had a reference result to compare our results with. This addressed the second problem, as, whatever the lists, we were confronted with the same challenge as the resource developers. We found datasets for seven languages and six domains. In each case we entered the corpus into the Sketch Engine, ran the term-finder, and computed precision and recall (which we could then compare with the performance figures of the group who developed the dataset.) The paper will present these results.

### 3. In sum

The Sketch Engine has for some years been a leading tool for lexicography and corpus linguistics. Its terminology function is now a year old. We present a thorough evaluation.

### References

- M. Baroni and S. Bernardini. 2004. [BootCaT: Bootstrapping corpora and terms from the web](#). Proceedings of LREC 2004, Lisbon: ELDA. 1313-1316.
- S. Bernardini, A. Ferraresi and E. Zanchetta. 2013. Old needs, new solutions: comparable corpora for language professionals. In Sharoff, S., R. Rapp, P. Zweigenbaum, P. Fung, editors. Building and Using Comparable Corpora. Springer